



# Bad Words: Finding Faults in Spirit's Syslogs

Resilience08 Workshop  
CCGrid08, Lyon France  
May 22, 2008

Jon Stearley  
[\*jrstear@sandia.gov\*](mailto:jrstear@sandia.gov)  
Sandia National Laboratories (US)



# Production Impacts

---

**Sisyphus has found:**

**Malfunctions:**

disks, controllers, network interfaces, power supplies, memory

**Misuse:**

RAID stripe imbalance, inappropriate remote monitoring

**Misconfigurations:**

BIOS, RAID controller, inconsistent software versions, config typos

**Which has enabled focused reactive and proactive responses.**

**Deployments:**

**SNL:** Red Storm, Thunderbird, Spirit, *TLCC, Corporate IT*

**LANL** *[monitoring suite]: TLCC, Roadrunner*

**450 Downloads** (as of 5/5/08)

See <http://www.cs.sandia.gov/sisyphus> for more info.

```
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Medium Error Disk 4G 3KT1HVCG Key: 3 ASC 16 ASCQ 0 FRU D2 Sense 80008E Info 0889A800
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889a800 LUN 7, 00000090 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889aa00 LUN 7, 00000091 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889ac00 LUN 7, 00000092 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889ac00 LUN 7, 00000093 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b000 LUN 7, 00000094 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b200 LUN 7, 00000095 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b400 LUN 7, 00000096 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 889b600 LUN 7, 00000097 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000090 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000091 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000092 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000093 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000094 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:02 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 7, 00000095 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
```

**Syslogs are:**

**Ubiquitous! Informational! Repetitive! Vast!**

**But how do you find the few lines of key  
information among thousands of log files and  
millions of lines of time-stamped text???**

```
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245800 LUN 6, 0001122b DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245800 LUN 6, 0001122c DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245a00 LUN 6, 0001122d DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245c00 LUN 6, 0001122e DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2245e00 LUN 6, 0001122f DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2246000 LUN 6, 00011230 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2246200 LUN 6, 00011231 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info INT_DG Data recovered disk:4G address: 2246400 LUN 6, 00011232 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r1 w0 l0 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122b DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122c DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122d DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122e DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 0001122f DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011230 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011231 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
Jul 16 10:00:07 10.1.0.49 local7 info DMT_EMT EMT verify reassign 1: LUN 6, 00011232 DLR:0, DLG:0, DRR:0, DEL:0, DELR:0, DERR:0 r0 w0 l1 fl0 fr2 ea:0,10
```



# Anomaly Detection in System Logs

---

## **Goal:**

Automatically detect “alerts” in system logs  
(messages of interest, eg malfunction or misuse).

## **Approach:**

Similar computers correctly executing similar work  
should produce similar logs  
(anomalies are “interesting”).

## **Measure:**

Quantify detection performance, using known  
signatures (regular expressions) as ground truth.



# Nodeinfo Algorithm

---

1. Group messages from N nodes over H hours into NH nodehour “docs” (docs/YYYY/MM/DD/HH/NODE)

2. Index to form term-doc matrix X  
(M terms by NH nodehours)

$term_i = \text{“PositionWord”}$   
e.g. “0003error”

3. Form term-node index Y (M terms by N nodes)

4. Using Y, calculate term information weights G  
(M by M diagonal)

$$g_i = 1 + \frac{1}{\log_2 N} \sum_{j=1}^N p_{ij} \log_2(p_{ij})$$
$$p_{ij} = \frac{y_{ij}}{\sum_{j=1}^N y_{ij}}$$

5. Rank docs by column magnitudes of  $G \log_2(X)$

$$nodeinfo_j = \sqrt{\sum_{i=1}^M (g_i \log_2(x_{ij}))^2}$$

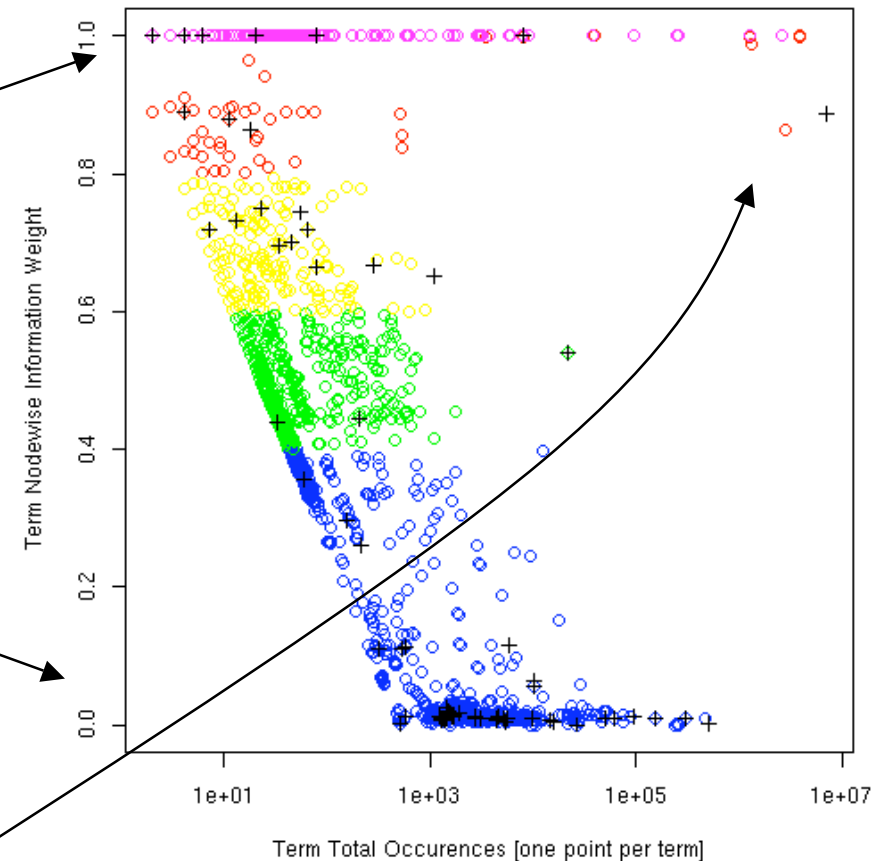


# Term Information Weights G

$g_i = 1$   
if term  $i$  occurs on only  
one node

$g_i = 0$   
if term  $i$  is distributed  
equally across all nodes

High-information terms  
occurring many times  
are most significant.



$$g_i = 1 + \frac{1}{\log_2 N} \sum_{j=1}^N p_{ij} \log_2(p_{ij})$$

$$p_{ij} = \frac{y_{ij}}{\sum_{j=1}^N y_{ij}}$$



# 0003error vs 0006error

```
May 20 23:01:00 sn105/sn105 CROND[*]: LAuS error - do_command.c:226 - laus_attach: (19) laus_attach: No such device
May 20 23:46:42 sn105/sn105 kernel: EXT3-fs error (device cciss0(104,2)): ext3_get_inode_loc: unable to read inode block
May 20 23:46:42 sn105/sn105 kernel: EXT3-fs error (device cciss0(104,2)): ext3_get_inode_loc: unable to read inode block
May 20 23:46:42 sn105/sn105 Event Log Daemon:[2907]: Fatal drive error, SCSI port 1 ID 0
May 20 23:46:43 sn105/sn105 Event Log Daemon:[2907]: Fatal drive error, SCSI port 1 ID 0
```

<u>pos</u>	<u>word</u>	<u>host info</u>	<u>count</u>	<u>support</u>	<u>host weight</u>	<u>host count</u>
9	<a href="#">laus_attach:</a>	0.00	1	248058	0.001628	508
10	<a href="#">No</a>	0.00	1	248058	0.001628	508
11	<a href="#">such</a>	0.00	1	248058	0.001628	508
12	<a href="#">device</a>	0.00	1	248058	0.001628	508
3	<a href="#">error</a>	9.44	1934	2810440	0.864250	508
4	<a href="#">Fatal</a>	3.91	15	15	1.000000	1
6	<a href="#">error.</a>	3.91	15	15	1.000000	1
7	<a href="#">unable</a>	9.92	967	1281749	1.000000	1
9	<a href="#">read</a>	9.92	967	1281749	1.000000	1
10	<a href="#">inode</a>	9.92	967	1281749	1.000000	1
11	<a href="#">block</a>	9.92	967	1281749	1.000000	1

**Not always an alert!**

**Always an alert!**

In above  
messages

Across all  
docs

Out of  
512 hosts





Reboots cause bursts of messages, most of which are not important.

But in this case, there was an inconsistent BIOS setting!

“0001kernel:” is occurs in many alerts, and many non-alerts (and contributes to false alarms if not ignored).

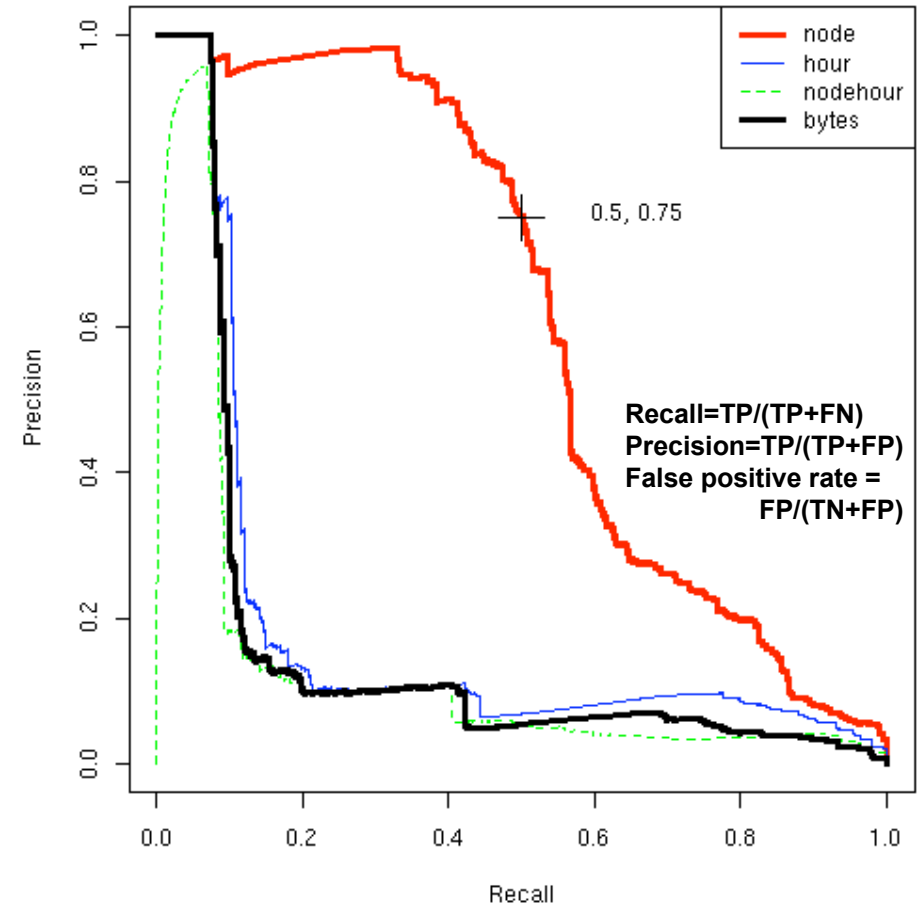
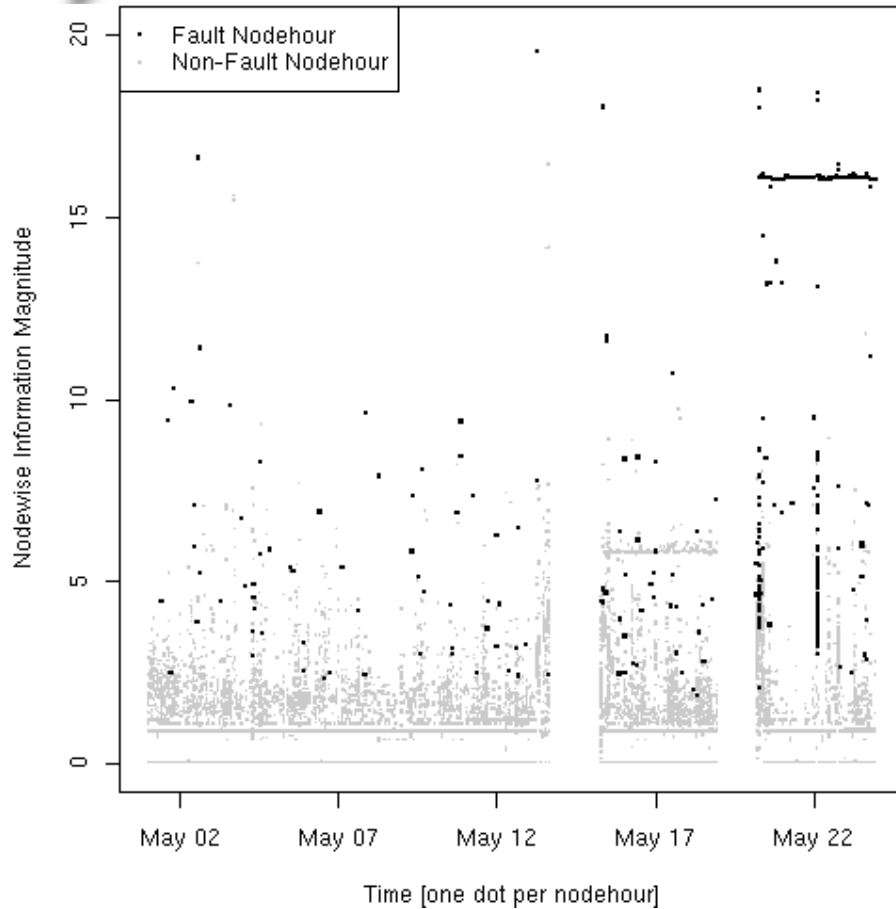
```
kernel: CPU1: Intel(R) Xeon(TM) CPU 3.40GHz stepping 04
kernel: Booting processor 2/6 rip 6000 page 0000010008764000
kernel: Initializing CPU#2
kernel: Calibrating delay loop... 6763.31 BogoMIPS
kernel: Monitor/Mwait feature present.
kernel: CPU: Trace cache: 12K uops<6>CPU: L2 cache: 1024K
kernel: CPU: Physical Processor ID: 3
kernel: Intel machine check reporting enabled on CPU#2.
kernel: CPU2: Intel(R) Xeon(TM) CPU 3.40GHz stepping 04
kernel: Booting processor 3/7 rip 6000 page 000001007ffea000
kernel: Initializing CPU#3
kernel: Calibrating delay loop... 6789.52 BogoMIPS
kernel: Monitor/Mwait feature present.
kernel: CPU: Trace cache: 12K uops<6>CPU: L2 cache: 1024K
kernel: CPU: Physical Processor ID: 3
kernel: Intel machine check reporting enabled on CPU#3.
kernel: CPU3: Intel(R) Xeon(TM) CPU 3.40GHz stepping 04
kernel: Total of 4 processors activated (27131.90 BogoMIPS).
kernel: cpu_sibling_map[0] = 1
kernel: cpu_sibling_map[1] = 0
kernel: cpu_sibling_map[2] = 3
kernel: cpu_sibling_map[3] = 2
kernel: mapping CPU#0's runqueue to CPU#1's runqueue.
kernel: mapping CPU#2's runqueue to CPU#3's runqueue.
```

pos	word	host	info	count	support	host	weight	host	count	time	weight	time	count	doc	weight	doc	coun
1	kernel:	7.63		389	6966438	0.886604		512		0.445072		247		0.692564		4147	
2	CPU1:	0.00		1	1286	0.009201		511		0.696248		40		0.424128		1272	
8	CPU#2.	0.00		1	2	1.000000		1		0.887916		2		0.944112		2	
2	CPU2:	0.00		1	2	1.000000		1		0.887916		2		0.944112		2	





# Nodehour Information Magnitudes



**Nodeinfo outperforms bytes.**  
**Hourinfo and Docinfo do not.**  
Nor does tf.idf weighting (not shown).

$$nodeinfo_j = \sqrt{\sum_{i=1}^M (g_i \log_2(x_{ij}))^2}$$



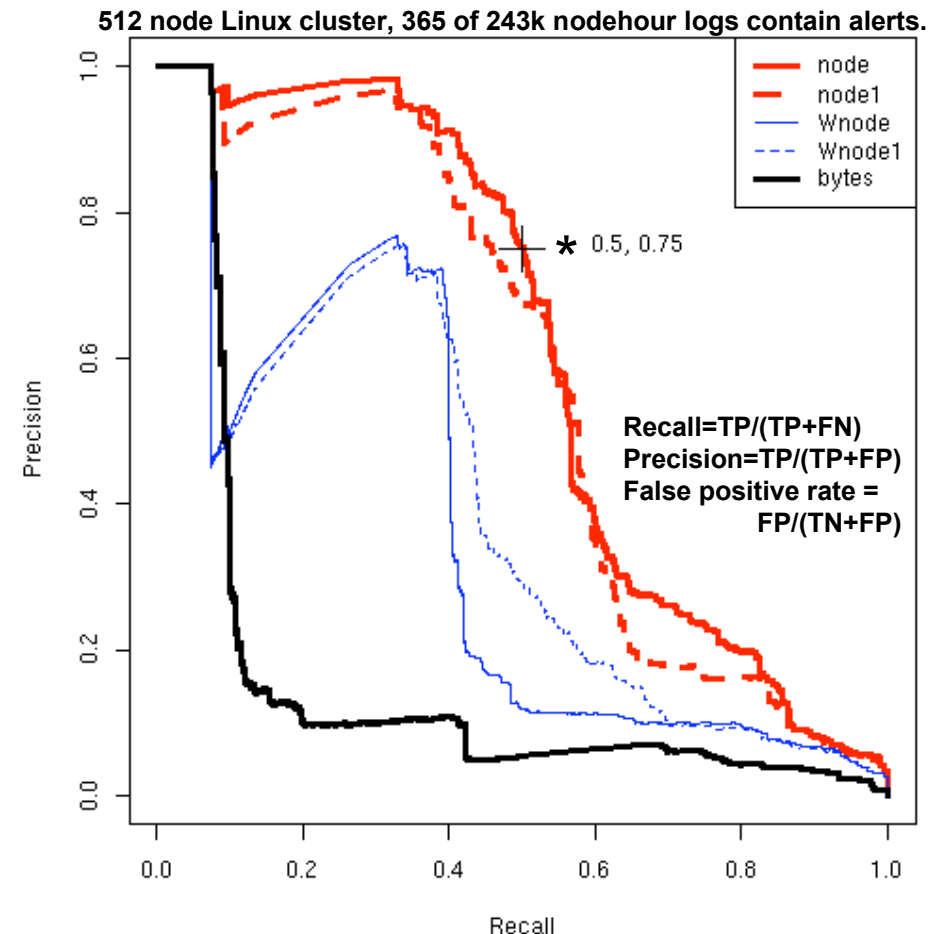
# Conclusions

**Bytes only detects message bursts (alerts, or not).**

**Nodeinfo detects more types of alerts.\***

**Word position information is significant.**  
(terms vs words)

**Ignore first words (dashed).**  
(set  $g_i=0$  for “0001” terms)



**\* 75% precision at 50% recall, corresponding to an excellent false-positive rate of 0.05%.**



# Open Questions

---

**Would a combination of nodeinfo and timeinfo and docinfo would be more effective?**

**Are we destroying too much context by capturing only word position information?**  
(e.g. explore term n-grams or message n-grams?)

**Terms are regular expressions (RE's) plus position information - what a pain to use and tune!**

- Are terms too burdensome in practice?**
- Are RE's rich enough to describe all anomalies of interest?**

**E.g. how to *predict* them *before* they occur???**



## Take Aways

---

**Nodeinfo is computationally simple and effective at detecting a wide range of alert messages.**

**Sisyphus is used on production supercomputers at SNL, and is publicly downloadable (LGPL) at <http://www.cs.sandia.gov/sisyphus>.**

**Logs are a rich mountain to mine for resilience!**



# Extra slides follow...

---

Sisyphus – Document Selection

https://rssweb.sandia.gov/ddn/docs.html

Sisyphus – Document Selection Sisyphus – Document Analysis Sisyphus – Document Analysis

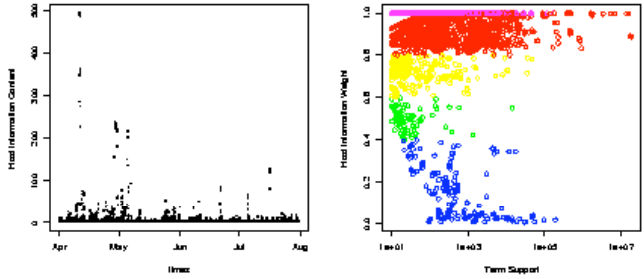
ddn goto destination

Date: start: 2007/07/01 10:00:00 First: 2007/04/01 00:00:00  
stop: 2007/07/31 11:00:00 Last: 2007/07/31 11:00:00  
Host(s): (not ☐)  
Word(s): (or ☐)  
Display top N: 100 sorted by host\_info Desc.

List: Documents ?

Displaying 100 of 42908 docs. 0 selected. (0 Kbytes) [HELP] [change settings]

Analyze: Templates Msgs. Terms ?



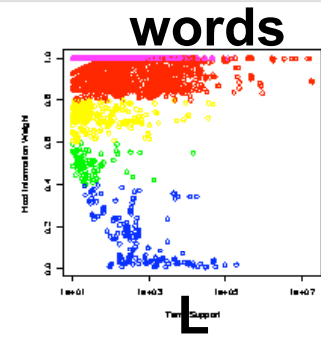
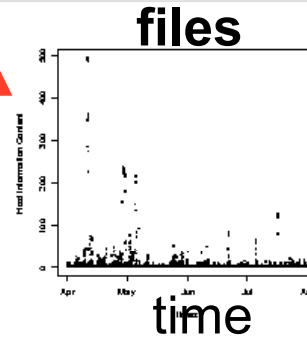
# 1. Which log files contain useful information?



ddn goto destination

Date: start: 2007/07/01 10:00:00 First: 2007/04/01 00:00:00  
 stop: 2007/07/31 11:00:00 Last: 2007/07/31 11:00:00  
 Host(s): (not ☐)  
 Word(s): (or ☐)  
 Display top N: 100 sorted by host\_info Desc.

$|(\mathbf{GL})_j|$



G

List: Documents ?

Displaying 100 of 42908 docs. 0 selected. (0 Kbytes) [HELP] [change settings]

Analyze: Templates Msgs. Terms ?

	YYYY/MM/DD/HH	HOST	bytes	lines	host_info	time_info	doc_info
<input type="checkbox"/>	docs/2007/07/16/10/10.1.0.49		1842113	11251	130.493	102.764	111.610
<input type="checkbox"/>	docs/2007/07/16/09/10.1.0.49		1867390	11437	129.803	102.591	111.241
<input type="checkbox"/>	docs/2007/07/16/11/10.1.0.49		1816549	11125	129.538	102.359	111.008
<input type="checkbox"/>	docs/2007/07/16/12/10.1.0.49		1704339	10437	126.824	100.048	108.612
<input type="checkbox"/>	docs/2007/07/16/08/10.1.0.49		1481224	9068	120.769	95.549	103.661
<input type="checkbox"/>	docs/2007/07/16/13/10.1.0.49		430320				
<input type="checkbox"/>	docs/2007/07/05/09/10.1.0.12		288005				
<input type="checkbox"/>	docs/2007/07/05/09/10.1.0.16		158502				
<input type="checkbox"/>	docs/2007/07/05/09/10.1.0.11		77539				
<input type="checkbox"/>	docs/2007/07/05/09/10.1.0.6		17907				
<input type="checkbox"/>	docs/2007/07/22/14/10.1.0.35		4430				
<input type="checkbox"/>	docs/2007/07/22/05/10.1.0.35		4210				
<input type="checkbox"/>	docs/2007/07/27/12/10.1.0.47		4887				
<input type="checkbox"/>	docs/2007/07/22/13/10.1.0.35		3324				
<input type="checkbox"/>	docs/2007/07/05/10/10.1.0.28		11386				
<input type="checkbox"/>	docs/2007/07/29/19/10.1.0.47		3746				
<input type="checkbox"/>	docs/2007/07/05/12/10.1.0.38		11966				
<input type="checkbox"/>	docs/2007/07/27/14/10.1.0.47		3526				
<input type="checkbox"/>	docs/2007/07/23/17/10.1.0.35		2660				
<input type="checkbox"/>	docs/2007/07/05/12/10.1.0.46		10482				
<input type="checkbox"/>	docs/2007/07/05/12/10.1.0.51		10653				
<input type="checkbox"/>	docs/2007/07/05/09/10.1.0.7		9918				
<input type="checkbox"/>	docs/2007/07/05/14/10.1.0.46		9172				
<input type="checkbox"/>	docs/2007/07/05/14/10.1.0.38		9637				
<input type="checkbox"/>	docs/2007/07/05/14/10.1.0.51		9480				
<input type="checkbox"/>	docs/2007/07/24/10/10.1.0.35		2218				
<input type="checkbox"/>	docs/2007/07/05/13/10.1.0.46		6495				
<input type="checkbox"/>	docs/2007/07/05/12/10.1.0.7		6688				
<input type="checkbox"/>	docs/2007/07/05/13/10.1.0.28		6440				
<input type="checkbox"/>	docs/2007/07/05/12/10.1.0.28		6130				

abnormal

“interestingness”

normal

“interestingness”

(aka “information”) is purely mathematical ( $=|(\mathbf{GL})_j|$ ).

$$G_{i,j} = 1 + H_i, \quad L = \log_2(\mathbf{tf}_{i,j})$$

$$H_i = \sum_j p_{ij} \log_2(p_{ij}) / \log_2(d)$$

where  $p_{ij} = \mathbf{tf}_{i,j} / \sum_j \mathbf{tf}_{i,j}$

and  $\mathbf{tf}_{i,j}$  is how many times the  $i$ 'th word occurs in the  $j$ 'th file







Sisyphus - Document Analysis

https://rssweb.sandia.gov/ddn/analyze.cgi?file=docs/2007/07/05/14/10.1.0.46&analyze=messages

Sisyphus - Document Selection Sisyphus - Document Analysis Sisyphus - Document Analysis

<docs/2007/07/05/14/10.1.0.46> Email URL 3 templates, minsup=20, 49 terms, 0 selected. [change settings] Analyze: Templates Msgs. Terms

```

Jul 5 14:51:11 10.1.0.46 local7 info BIT_MON Host port 2 : SFP signal OK
Jul 5 14:51:13 10.1.0.46 local7 info BIT_MON Host port 3 : SFP signal OK
Jul 5 14:51:13 10.1.0.46 local7 info BIT_MON Host port 4 : SFP signal OK
Jul 5 14:51:17 10.1.0.46 local7 info VER_14 Verify disk errors on LUN 14 due to disk 8F error.
Jul 5 14:52:03 10.1.0.46 local7 info BIT_MON Host port 3 : SFP loss of signal
Jul 5 14:52:17 10.1.0.46 local7 info VER_15 Verify disk errors on LUN 15 due to disk 8F error.
Jul 5 14:53:12 10.1.0.46 local7 info TEL_MAIN New TELNET Session initiated from IP address: 10.1.0.121
Jul 5 14:53:13 10.1.0.46 local7 info tTelnetI -- remote Telnet session disconnect --
Jul 5 14:53:13 10.1.0.46 local7 info TEL_EXIT Telnet Session termination.
Jul 5 14:53:17 10.1.0.46 local7 info VER_15 Verify disk errors on LUN 15 due to disk 8F error.
Jul 5 14:53:57 10.1.0.46 local7 info TEL_MAIN New TELNET Session initiated from IP address: 10.1.0.121
Jul 5 14:53:58 10.1.0.46 local7 info tTelnetI -- remote Telnet session disconnect --
Jul 5 14:53:58 10.1.0.46 local7 info TEL_EXIT Telnet Session termination.
Jul 5 14:54:12 10.1.0.46 local7 info TEL_MAIN New TELNET Session initiated from IP address: 10.1.0.121
Jul 5 14:54:12 10.1.0.46 local7 info tTelnetI -- remote Telnet session disconnect --
Jul 5 14:54:12 10.1.0.46 local7 info TEL_EXIT Telnet Session termination.
Jul 5 14:54:17 10.1.0.46 local7 info VER_14 Verify disk errors on LUN 14 due to disk 8F error.
Jul 5 14:54:28 10.1.0.46 local7 info TEL_MAIN New TELNET Session initiated from IP address: 10.1.0.121
Jul 5 14:54:28 10.1.0.46 local7 info tTelnetI -- remote Telnet session disconnect --
Jul 5 14:54:28 10.1.0.46 local7 info TEL_EXIT Telnet Session termination.
Jul 5 14:55:17 10.1.0.46 local7 info VER_14 Verify disk errors on LUN 14 due to disk 8F error.
Jul 5 14:56:17 10.1.0.46 local7 info VER_14 Verify disk errors on LUN 14 due to disk 8F error.
Jul 5 14:56:25 10.1.0.46 local7 info TEL_MAIN New TELNET Session initiated from IP address: 10.1.0.121
Jul 5 14:56:25 10.1.0.46 local7 info tTelnetI -- remote Telnet session disconnect --
Jul 5 14:56:25 10.1.0.46 local7 info TEL_EXIT Telnet Session termination.
Jul 5 14:57:03 10.1.0.46 local7 info TEL_MAIN New TELNET Session initiated from IP address: 10.1.0.121

```

<input type="checkbox"/>	<input checked="" type="checkbox"/>	ID	pos	word	host info	count	support	host weight	host count	time weight	time count	doc weight	doc count
<input type="checkbox"/>	<input checked="" type="checkbox"/>	2	2	info	5.92	99	16944740	0.89232	98	0.6413	2547	0.7265	42864
<input type="checkbox"/>	<input type="checkbox"/>	13	13	8F	5.91	60	150	1.00000	1	0.8540	4	0.8927	4
<input type="checkbox"/>	<input type="checkbox"/>	8	8	LUN	4.93	60	901319	0.83446	27	0.4870	25	0.6225	386
<input type="checkbox"/>	<input type="checkbox"/>	9	9	15	4.91	30	72			851			4
<input type="checkbox"/>	<input type="checkbox"/>	16	9	14	4.90	30	182608			540			49
<input type="checkbox"/>	<input type="checkbox"/>	6	6	errors	3.80	60	947			785			24
<input type="checkbox"/>	<input type="checkbox"/>	10	10	due	3.80	60	947			785			24
<input type="checkbox"/>	<input type="checkbox"/>	14	14	error.	3.80	60	947	0.64265	7	0.7854	9	0.7232	24
<input type="checkbox"/>	<input type="checkbox"/>	12	12	disk	3.77	60	951	0.63746	9	0.7825	11	0.7216	26
<input type="checkbox"/>	<input type="checkbox"/>	11	11	to	3.65	60	974	0.61950	11	0.7717	12	0.7162	20

on 1 computer (out of 90)

over 4 hours (out of 4 months)

Useful term statistics.



# Useful Patterns

## Automatically generated message templates and time statistics.

<input type="checkbox"/>	<a href="#">ID</a>	<a href="#">count</a>	<a href="#">median</a>	<a href="#">stddev</a>	<a href="#">regex</a>
<input checked="" type="checkbox"/>	0	113	0	57	<a href="#">OUTLIERS</a>
<input type="checkbox"/>	1	40	15	100	daemon info llrd 5640 : llrd: nid00192 - - 17/Oct/2007 * "POST /RPC2 HTTP/1.0" 200 -
<input type="checkbox"/>	2	20	0	389	kern * kernel: * * slow * *
<input type="checkbox"/>	3	9	301	472	kern err kernel: LustreError: * * * *
<input type="checkbox"/>	4	47	13	102	kern alert kernel: LustreError: dumping log to *
<input type="checkbox"/>	6	4	760	853	kern * kernel: * dumping log to *
<input type="checkbox"/>	7	6	602	16	kern * kernel: * * * * *
<input type="checkbox"/>	8	86	22	132	kern warning kernel: SCSI error : <1 0 0 0> return code = 0x20000
<input type="checkbox"/>	18	86	22	132	kern warning kernel: end_request: I/O error, dev sde, sector *
<input type="checkbox"/>	20	50	0	97	kern warning kernel: Call Trace:{schedule_timeout+243} {process_timeout+0}
<input type="checkbox"/>	21	36	0	79	kern warning kernel: Call Trace:{libcfs:libcfs_nid2str+178} {ost:ost_brw_write+2000}
<input type="checkbox"/>	22	2	301	0	kern warning kernel: Call Trace:{libcfs:libcfs_nid2str+178} *
<input type="checkbox"/>	23	2	600	0	kern warning kernel: Call * {ost:ost_brw_write+2000}

```
Oct 17 05:04:06 nid00187 kern crit kernel: LDISKFS-fs error (device sde2) in ldiskfs_setattr: Readonly filesystem
Oct 17 05:04:12 nid00187 kern warning kernel: SCSI error : <1 0 0 0> return code = 0x20000
Oct 17 05:04:12 nid00187 kern warning kernel: end_request: I/O error, dev sde, sector 778694416
Oct 17 05:04:12 nid00187 kern err kernel: Buffer I/O error on device sde2, logical block 7372802
Oct 17 05:04:12 nid00187 kern warning kernel: lost page write due to I/O error on sde2
Oct 17 05:04:12 nid00187 kern warning kernel: SCSI error : <1 0 0 0> return code = 0x20000
Oct 17 05:04:12 nid00187 kern warning kernel: end_request: I/O error, dev sde, sector 779218704
Oct 17 05:04:12 nid00187 kern err kernel: Buffer I/O error on device sde2, logical block 7438338
Oct 17 05:04:12 nid00187 kern warning kernel: lost page write due to I/O error on sde2
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6388:0:(lustre_fsfilth:255:fsfilt_commit_wait()) slow journal start 51s
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(filter_io_26.c:707:filter_commitrw_write()) slow commitrw commit 3511s
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(filter_io_26.c:707:filter_commitrw_write()) previously skipped 5 similar messages
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(service.c:583:ptlrpc_server_handle_request()) request 527 ope 4 from U3-1251@ptl processed in 3511s trans 0
rc -5/-5
Oct 17 05:04:20 nid00187 kern err kernel: LustreError: 6388:0:(service.c:583:ptlrpc_server_handle_request()) previously skipped 7 similar messages
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6388:0:(watchdog.c:320:lcw_update_time()) Expired watchdog for pid 6388 disabled after 3511.0309s
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6339:0:(watchdog.c:320:lcw_update_time()) Expired watchdog for pid 6339 disabled after 3511.4820s
Oct 17 05:04:20 nid00187 kern warning kernel: Lustre: 6388:0:(watchdog.c:320:lcw_update_time()) previously skipped 7 similar messages
```





# Logs: Research Collaborations

---

## **Adam Oliner - Stanford**

Time and/or Space Correlated Anomalies

## **James Elliot, Box Leangsuksun - Louisiana Tech**

Latent Semantic Analysis

## **Risto Vaarandi - Cyberdefence Centre of Excellence (EU)**

Term Patterns

## **Within Sandia**

Graph Layout (VxOrd) - Shawn Martin

Corporate IT Security - Paiz, Parks, Sery





# HPC Resilience

---

**SNL momentum and support is increasing**  
(eg resilience was explicitly prioritized in '08 LDRD call).

**Scientific research, engineering, and operation  
requires standardized definitions and measurements.**

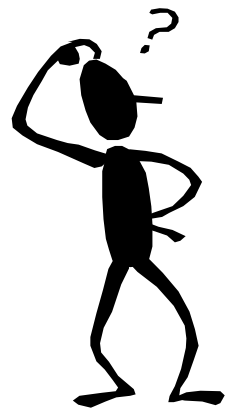
**Logs are a rich resilience research area.**  
Logs DO contain malfunction and misuse info.  
Current practices are painful and insufficient.



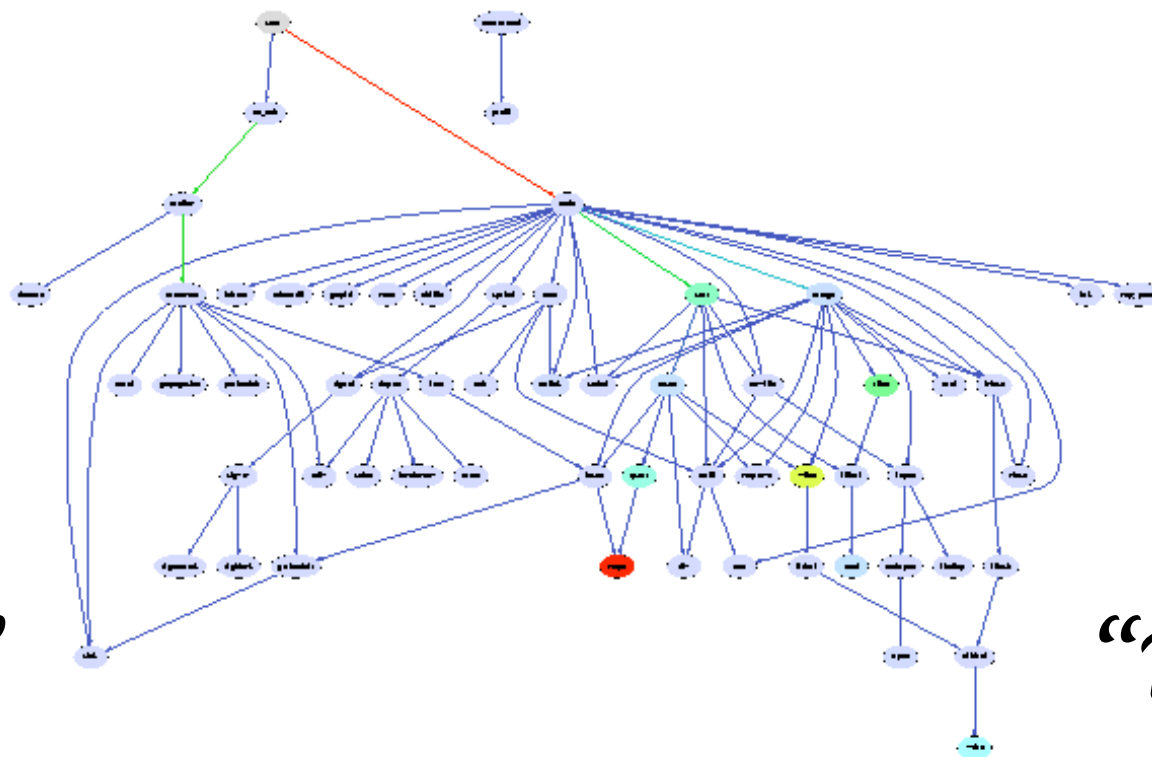
# Status Quo

**“A computer is in one of two situations. It is either known to be bad or it is in an unknown state.”**

**Mike Levine (PSC)**



*“Up!”*



*“Down!”*



# Standard Metrics: Needed

---

**Everyone uses the same terms (eg MTBF)  
but different definitions and measurements.**

- BAD PRACTICE!!! (eg procurements and operations)
- BAD SCIENCE!!! (eg quantify algorithm performance)

## **Challenges:**

1. Agree on definitions and measurements  
eg: from sysadmin, user, or manager perspective?
2. Change our spoken and written language.
3. Change necessary operational processes and procedures.



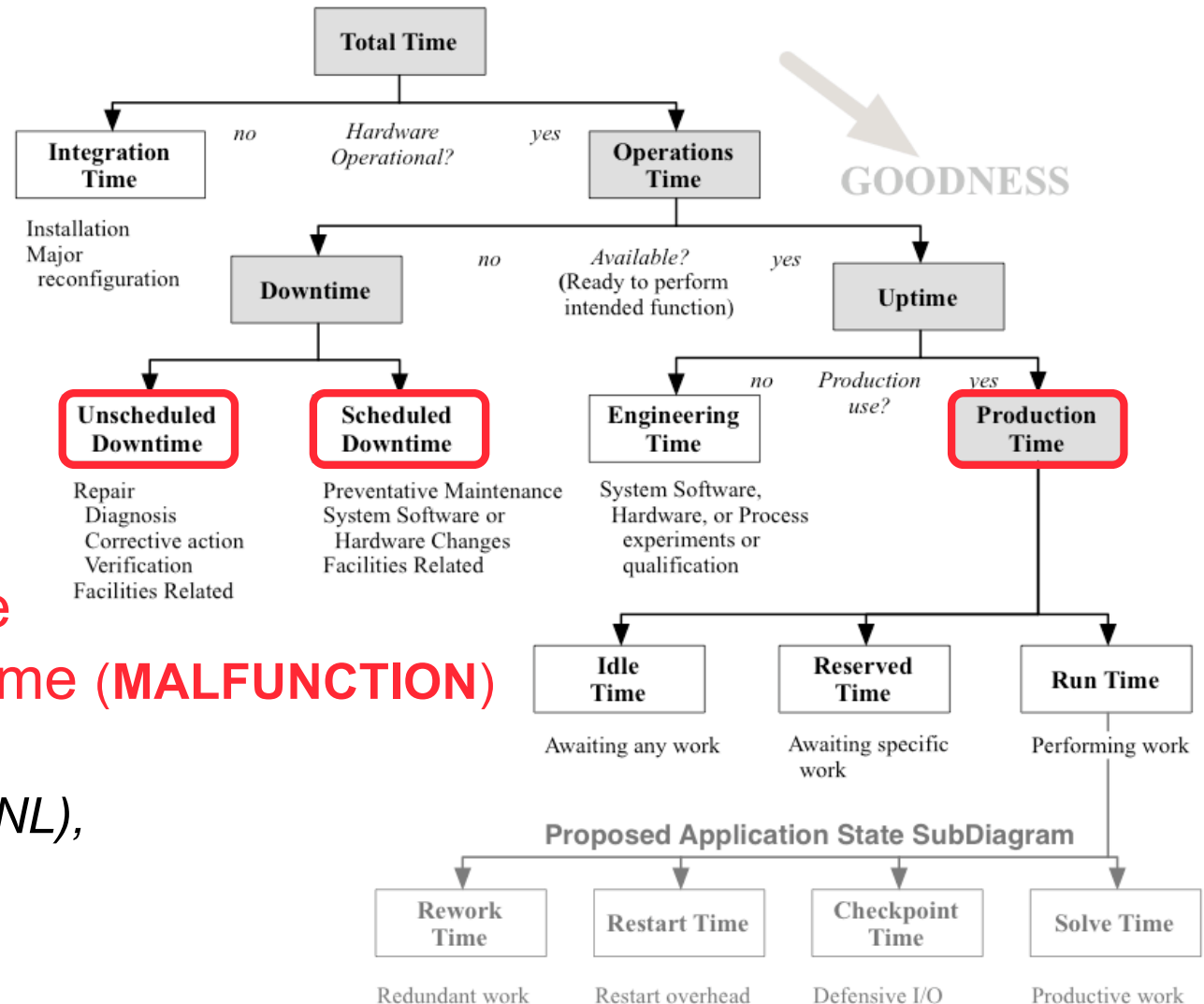
# Operations Status: Essential

**Tri-lab-developed Component State Diagram** (Based on SEMI-E10)  
Each component is in exactly one non-grey state at all times.

Need to log  
transitions of  
each node  
among three  
conditions:

Production Uptime  
Scheduled Downtime  
Unscheduled Downtime (**MALFUNCTION**)

Stearley (SNL), Daly (LANL),  
Hamilton (LLNL)





# Component Operations Status (COS)

---



Scheduled  
Downtime

Production  
Uptime

Unscheduled  
Downtime

## **Production Uptime (PU)**

ready for immediate use by one or more production user

## **Scheduled Downtime (SD)**

not in PU for scheduled reasons

## **Unscheduled Downtime (UD)**

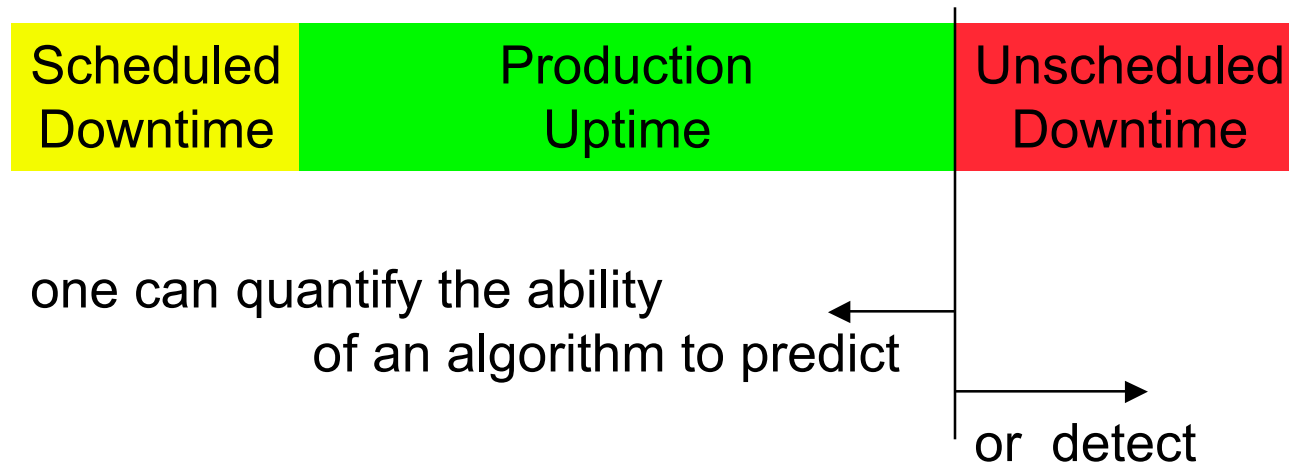
not in PU for unscheduled reasons



# Operations Status: Essential

---

Given per-component operations status data:



the onset of Unscheduled Downtime  
(by analyzing logs, or other data).